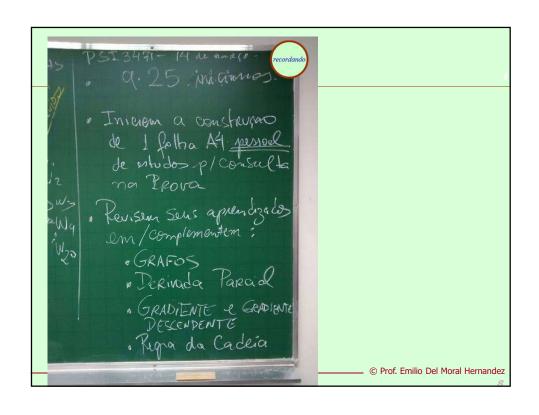
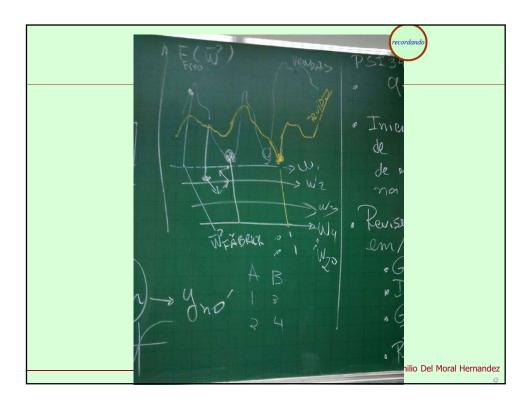
Aprendizado em RNAs do tipo MLP – Multi Layer Percetron – através do algoritmo Error Back Propagation

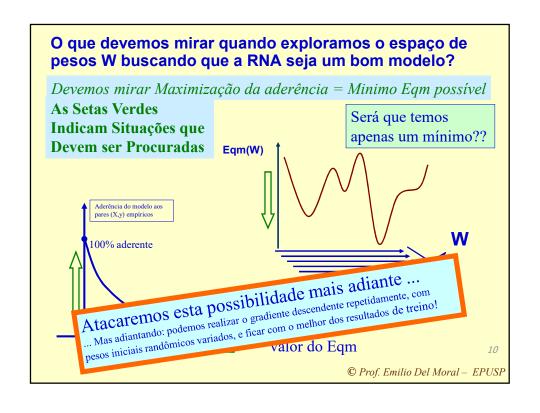
© Prof. Emilio Del Moral – EPUSP

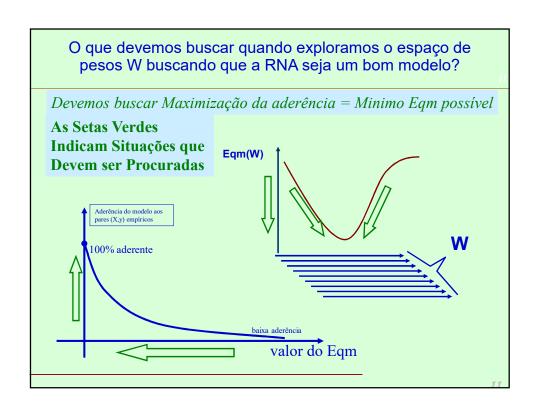
Como escolhemos os valores dos diversos w's?

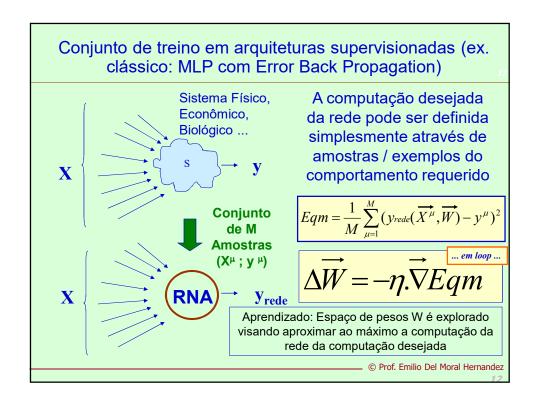
© Prof. Emilio Del Moral Hernandez

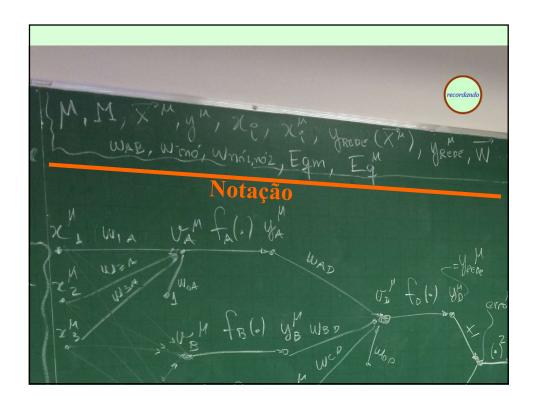


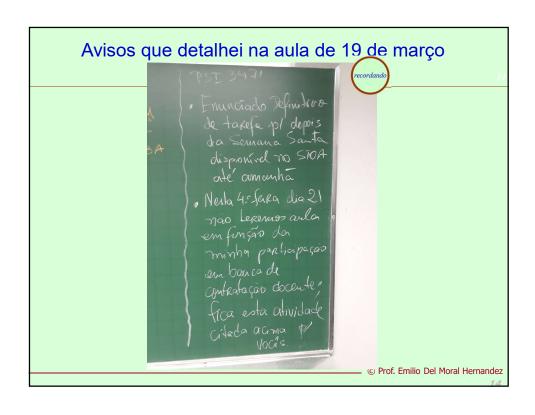


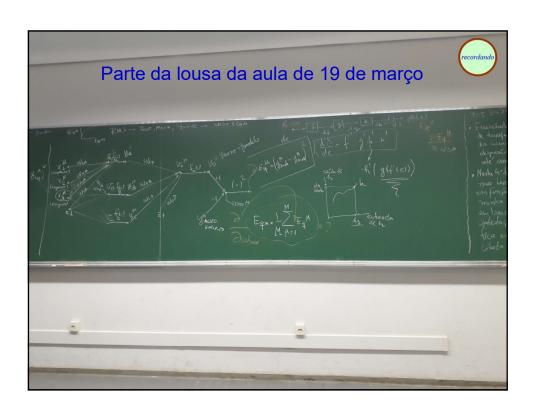




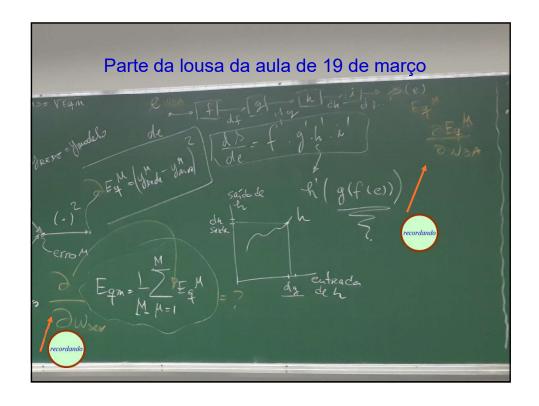








Deduzindo as Equações do Aprendizado em RNAs do tipo MLP – Multi Layer Percetron – com o algoritmo Error Back Propagation (Gradiente Descendente)



Enunciado c/ ligeira modificação em aula (para o seu melhor aprendizado)

Tarefa para desenvolvimento próprio e entrega até a aula de 4ª f. da Semana Santa +1(dia 4/abril)

- No grafo detalhado de uma rede 3-3-1 como o da aula, identifique o peso sináptico "w_{meu}" da primeira camada definido pelo último dígito do seu NUSP como segue:
- Último digito 1: w_{1A}; 2: w_{1B}; 3: w_{1C}; 4: w_{2A}; 5: w_{2B}; 6: w_{2C}; 7: w_{3A}; 8: w_{3B}; 9: w_{3C}; 0: w_{1A}
- Assumindo 100 exemplares de treino (M=100), identifique o exemplar de treino " μ_{meu} " associado ao 2º e 3º dígitos de seu NUSP como segue:
- Se 2° e 3° dígitos são 0 e 1, μ = 01, se 2° e 3° dígitos são 0 e 2, μ = 02, se 2° e 3° são 0 e 3, μ = 03 ... se 4 e 1, μ = 41 ... e assim por diante
- a) Fixando a entrada X da rede no exemplar empírico μ , ou seja X = X $^{\mu}$ e usando como alvo (target) empírico y = y $^{\mu}$, e visando o cálculo do gradiente do erro quadrático, necessário ao algoritmo de gradiente descendente (EBP error back propagation), deduza com detalhe a expressão analítica da derivada parcial do erro quadrático do "seu" exemplar μ_{meu} com relação ao "seu" peso sináptico w_{meu} , ou seja, calcule ($\partial Eq^{\mu_{meu}} / \partial w_{meu}$); revise seus conceitos de derivada parcial e use a regra da cadeia na sua dedução.
- b) Agora empregue a formula analítica obtida em a) para um (X^{μ_meu}, y^{μ_meu}) empírico com valores razoáveis ao seu sistema regressor (já criado por você nas atividades de aula): defina valores numéricos e unidades das 4 grandezas envolvidas no cálculo, x₁^μ, x₂^μ, x₃^μ, e y^μ, e chegue ao valor numérico de (∂Eq^{μ_meu} / ∂ w_{meu})
- c) Especifique o que teve que assumir em a e b (por não definido no enunciado).

© Prof. Emilio Del Moral Hernandez

O treinamento mira minimizar o Eqm das amostras (X; y) de treino. (exclusivamente!) sistema a modelar $\sum_{\mu} [y_{RNA}(X^{\mu}) - y_{sistema}^{\mu}]^2 / M$ Tabela de amostras Eam = Relação de amostragem ... (X; y) ... E lembremos que as amostras sempre são uma representação parcial do comportamento mais geral do sistema que está sendo modelado. (modelo do sistema) © Prof. Emilio Del Moral – EPUSP

Um Exemplo Ilustrativo para o Conceito de Conjunto de Treinamento e dos M pares (X,y)...

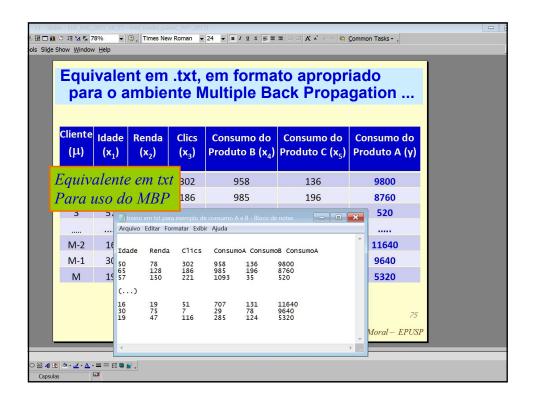
22

© Prof. Emilio Del Moral – EPUSP

Exemplo de regressão multivariada para estimação contínua usando MLP

- O valor do y contínuo ... neste exemplo corresponde ao volume de consumo futuro num dado tipo de produto "A" a ser ofertado pela empresa a um cliente corrente já consumidor de outros produtos da empresa ("B" e "C"), volume esse previsto com base em várias medidas quantitativas que caracterizam tal indivíduo. ... Assim, y = Consumo do Produto A = F(x₁,x₂, x₃, x₄, x₅).
- Consideremos 4 variáveis de entrada no modelo preditivo neural, ou seja, temos 5 medidas em X:
 - x₁: Idade do indivíduo
 - x₂: Renda mensal do indivíduo
 - x₃: Volume de clicks do indivído no website de exibição de produtos oferecidos pela empresa
 - x₄: Volume de consumo desse cliente observado para outro Produto B da mesma empresa
 - x₅: Volume de consumo desse cliente Produto C da mesma empresa
- Problema: desenvolver uma MLP para regressão contínua multivariada que permita estimar esse volume de consumo futuro y com base no conhecimento dos X e numa base de dados de aprendizado com esses dados X e y para 350 já clientes de universo populacional similar ao do novo consumidor potencial.

liente (µ)	Idade (x ₁)	Renda (x ₂)	Clics (x ₃)	Consumo do Produto B (x ₄)	Consumo do Produto C (x ₅)	Consumo do Produto A (y
1	50	78	302	958	136	9800
2	65	128	186	985	196	8760
3	57	150	221	1093	35	520
M-2	16	19	51	707	131	11640
M-1	30	75	7	29	78	9640
М	19	47	116	285	124	5320



A estratégia de Aprendizado para o MLP mais conhecida:

Error Back Propagation (EBP)

- = Propagação Reversa de Erro
- = Método do Gradiente personalizado ao Eqm(W) do MLP

30

© Prof. Emilio Del Moral – EPUSF

Mas entendamos PRIMEIRO
o que é o método numérico do
gradiente ascendente /
gradiente descendente
genérico,
que pode ser aplicado tanto para se
chegar paulatinamente ao máximo de
uma função quanto para se chegar ao
mínimo de uma função
(ascendente / descendente)

Chamada oral sobre a lição de casa: estudar / reestudar os conceitos e a parte operacional de derivadas parciais, do vetor Gradiente, e da regra da cadeia ...

 Derivadas parciais (que são as componentes do gradiente):

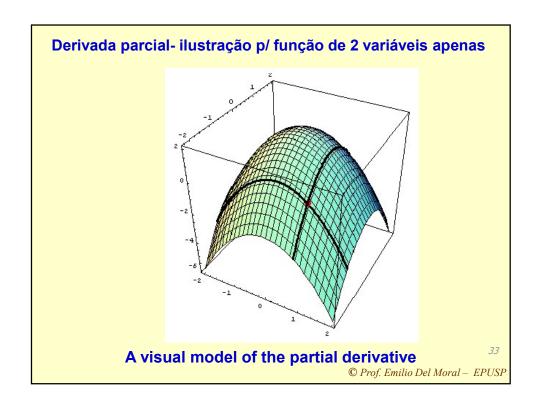
$$\partial f(a,b,c)/\partial a$$
 $\partial f(a,b,c)/\partial b$ $\partial f(a,b,c)/\partial c$

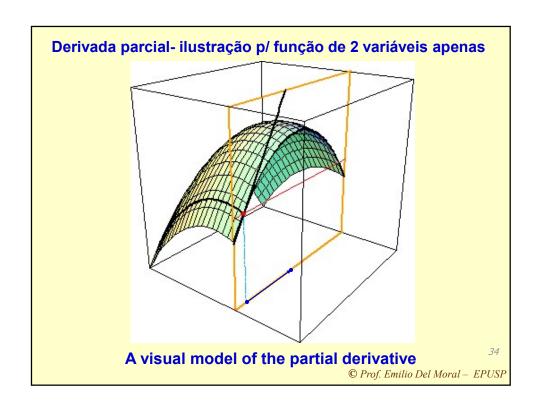
• Vetor Gradiente, útil ao método do máximo declive:

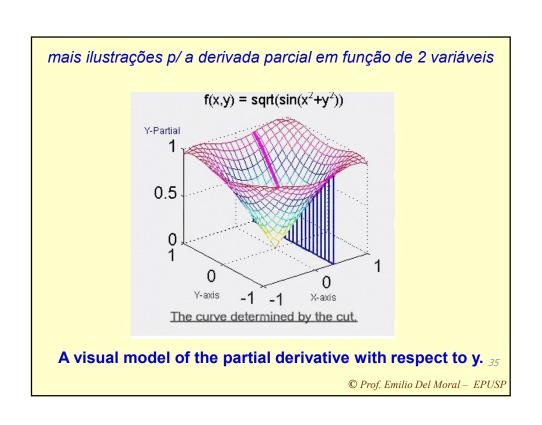
(
$$\partial \text{Eqm}(W)/\partial w_1$$
, $\partial \text{Eqm}(W)/\partial w_2$, $\partial \text{Eqm}(W)/\partial w_3$, ...)
$$\Delta \overrightarrow{W} = -\eta . \overrightarrow{\nabla} Eqm$$

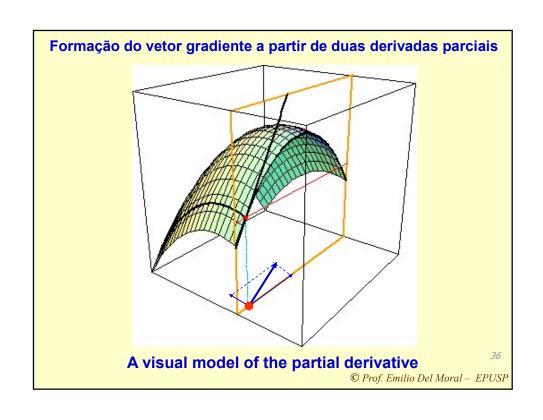
 Regra da cadeia, necessária ao cálculo de derivadas quando há encadeamento de funções:

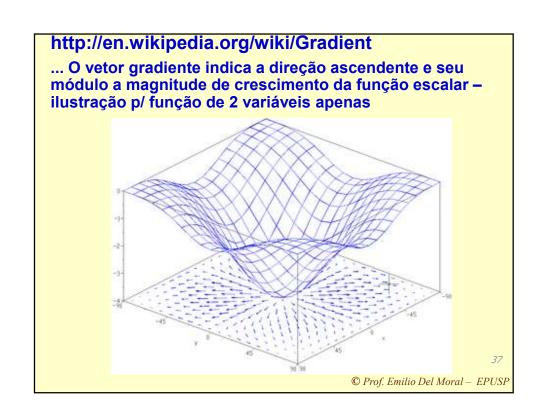
$$\partial f(g(h(a)))/\partial a = \partial f/\partial g \cdot \partial g/\partial h \cdot \partial h/\partial a$$





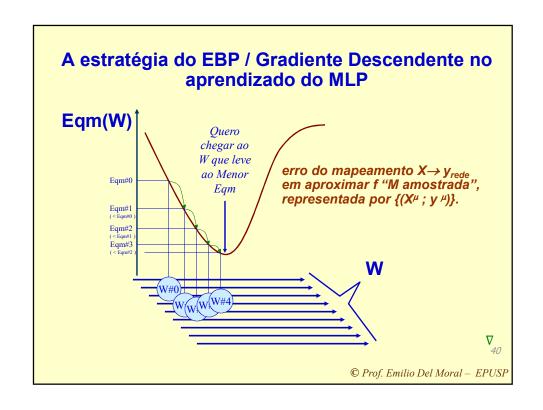






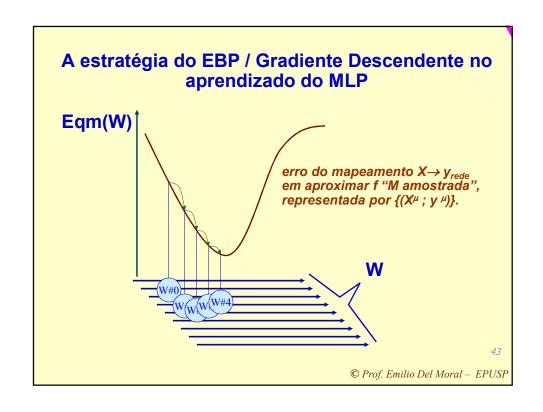


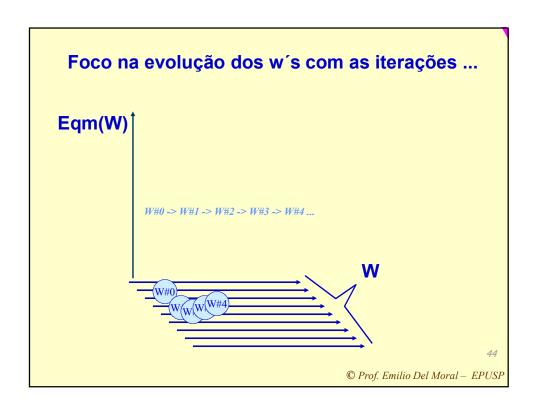
W#0	→ Eqm#0	GradEqm(W#0)	→ DeltaW#0 = - n.GradEqm(W#0
W#1 (= W#0 + DeltaW#0)	Eqm#1 (< Eqm#0)	GradEqm(W#1)	DeltaW#1 = - n.GradEqm(W#1
W#2 (= W#1 + DeltaW#1)	Eqm#2 (< Eqm#1)	GradEqm(W#2)	DeltaW#2 = - n.GradEqm(W#2
W#3 (= W#2 + DeltaW#2)	Eqm#3 (< Eqm#2)	GradEqm(W#3)	DeltaW#3 = - n.GradEqm(W#3
W#4 (= W#3 + DeltaW#3)	Eqm#4 (< Eqm#3)	GradEqm(W#4)	DeltaW#4 = - n.GradEqm(W#4
W#k (= W#k-1 + DeltaW#k-1)	Eqm#k (< Eqm#k-1)	GradEqm(W#k)	DeltaW#k = - n.GradEqm(W#4

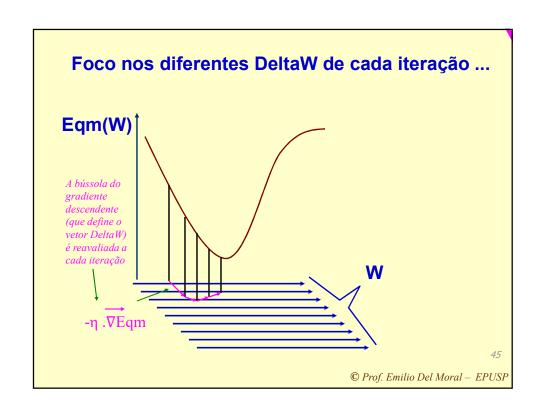


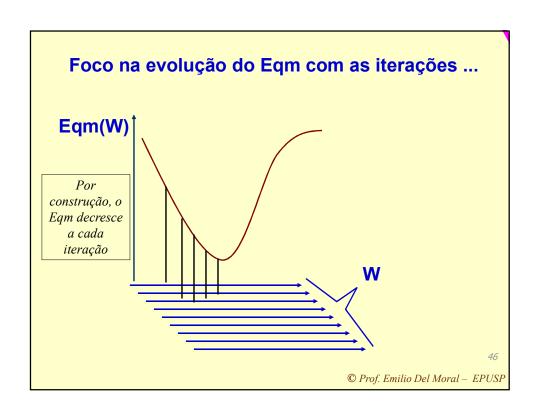
Este gráfico é denso e toca em muitos aspectos interrelacionados ... revisitemos alguns desses aspectos isoladamente com focos específicos nessas revisitas, assim teremos gráficos algo mais simples de interpretar ...

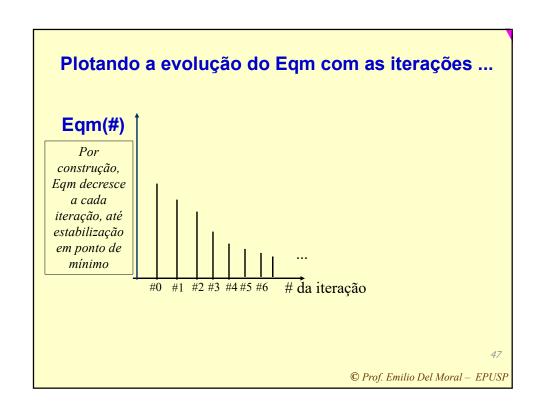
12

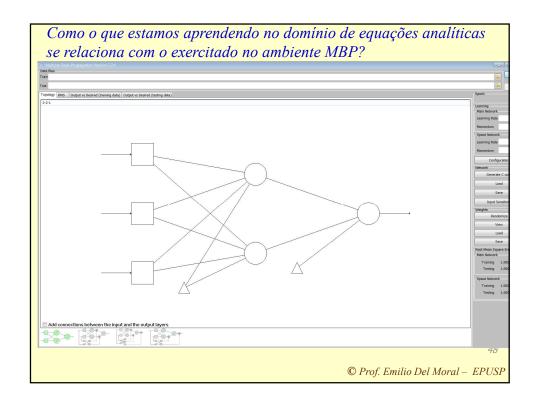


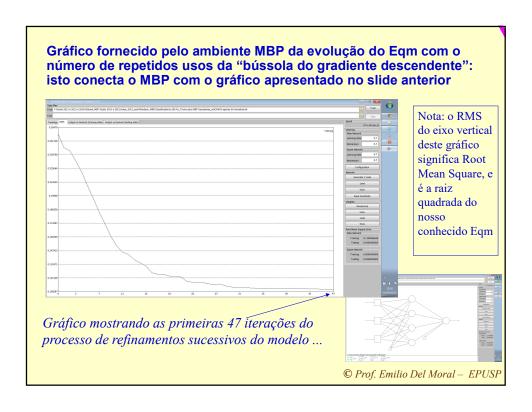


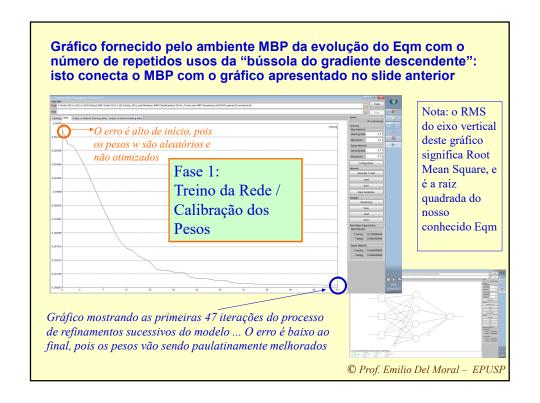












O Ciclo completo da modelagem:

- 0) Formalização do problema, mapeamento quantitativo em um modelo neural inicial e ... 0b) coleta de pares empíricos (X,v)
- 1) Fase de TREINO da RNA (MLP): com conhecimento dos X e dos y, que são ambos usados na calibração do modelo
- 2) Fase de TESTE / Caracterização da qualidade da RNA para generalizar: temos <u>novos</u> pares X e y, com y guardado "na gaveta", usado apenas para avaliação, não para re-calibração. É como um ensaio de uso final do modelo, com possibilidade de medir a sua qualidade com o y que foi guardado na gaveta.

[Fase de refinamentos da RNA, dados e modelo, em ciclos, desde 0]

3) Fase de USO FINAL da RNA, com y efetivamente não conhecido, e estimado com conhecimento dos X + uso do modelo calibrado.

.... Diferenças e semelhanças entre 1, 2 e 3

51

© Prof. Emilio Del Moral – EPUSP

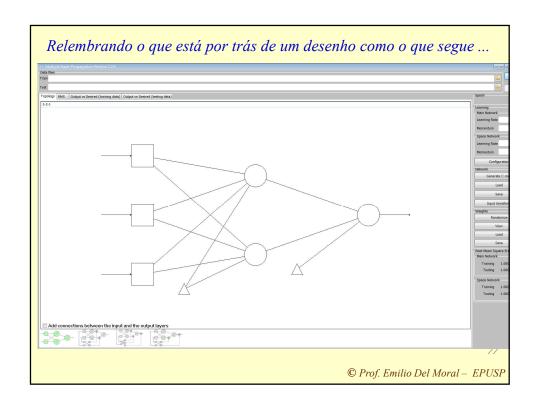
Aprendizado do MLP por Error Back Propagation ...

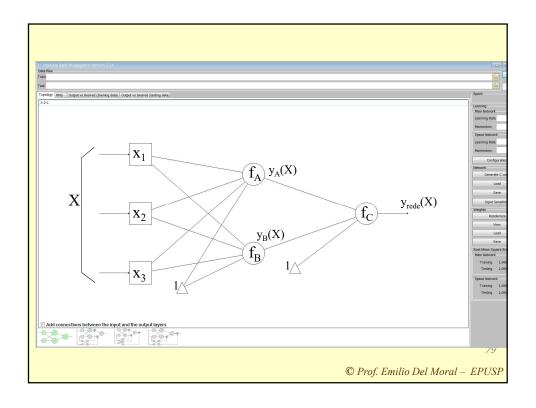
$$\Delta \overrightarrow{W} = -\eta . \overrightarrow{\nabla} Eqm_{\perp}$$

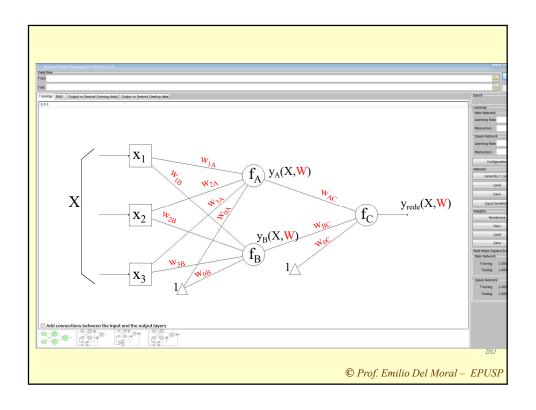
Gradiente de Eqm no espaço de pesos = $(\partial Eqm(W)/\partial w_1, \partial Eqm(W)/\partial w_2, \partial Eqm(W)/\partial w_3, ...)$

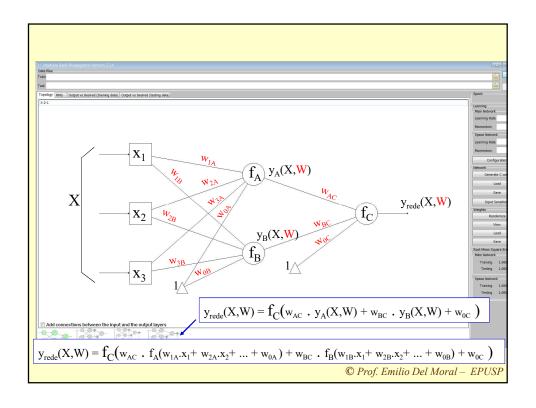
Chegando às fórmulas das derivadas parciais, necessárias à Bússola do Gradiente

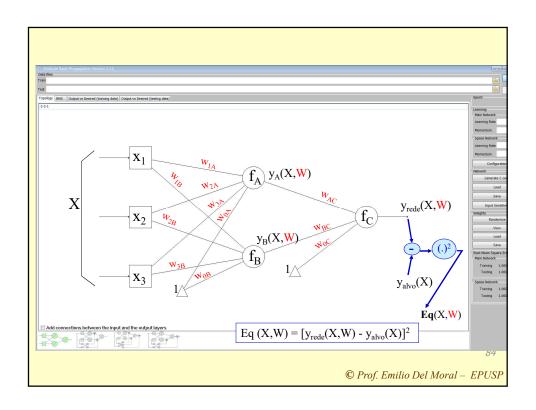
76

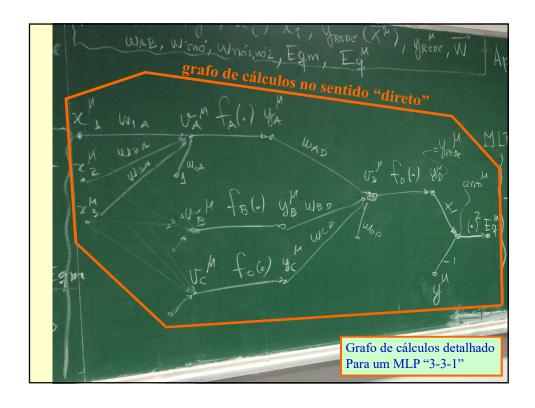












Chamada oral sobre a lição de casa: estudar / reestudar os conceitos e a parte operacional de derivadas parciais, do vetor Gradiente ...

Derivadas parciais (que são as componentes do gradiente):

 $\partial f(a,b,c)/\partial a$ $\partial f(a,b,c)/\partial b$ $\partial f(a,b,c)/\partial c$

• Vetor Gradiente, útil ao método do máximo declive:

 $(\partial Eqm(W)/\partial w_1, \partial Eqm(W)/\partial w_2, \partial Eqm(W)/\partial w_3, ...)$

$$\Delta \overrightarrow{W} = -\eta . \overrightarrow{\nabla} Eqm$$

86

© Prof. Emilio Del Moral – EPUSP

Invertamos o operador gradiente e a somatória

.. afinal, gradiente é uma derivada, e a derivada de um soma de várias funções é igual à soma das derivadas individuais de cada componente da soma:

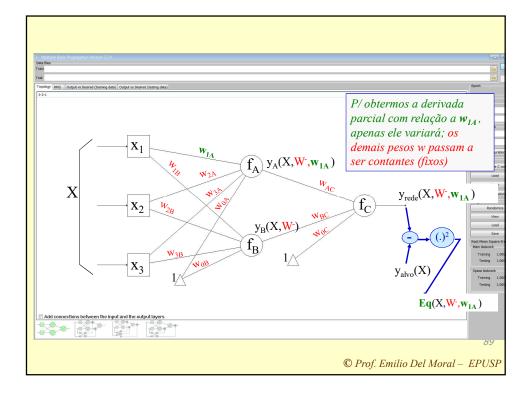
Grad(Eqm) =

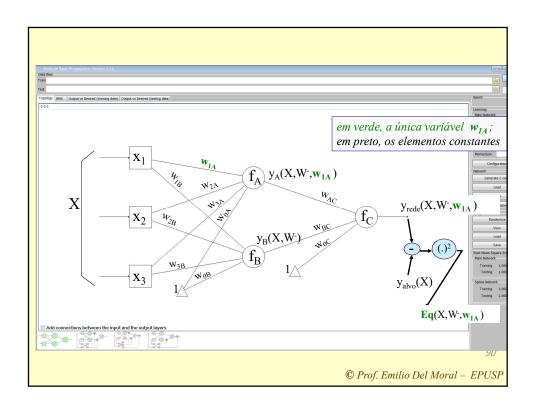
Grad($\Sigma_{\mu} Eq^{\mu}$) / M

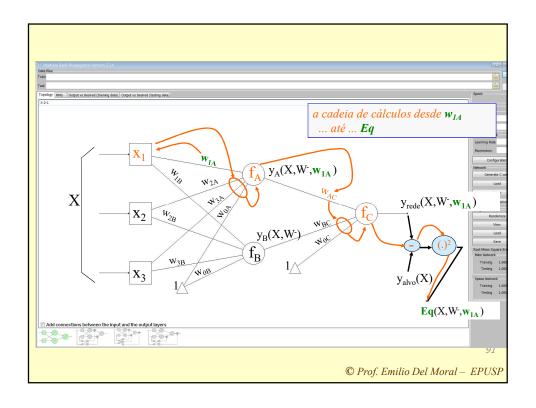
 Σ_{μ} Grad (Eq^{μ}) / M

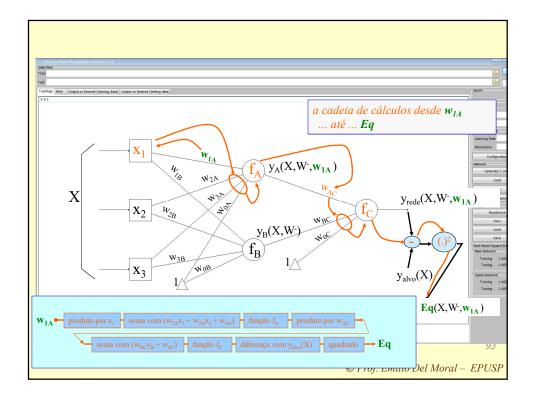
87

Note que a inversão do gradiente com a somatória nada mais é que usar de forma repetida – e em separado para cada dimensão do vetor $\textbf{Grad}(\Sigma_{\mu} \textbf{Eq}^{\mu})$ – a seguinte propriedade simples e sua velha conhecida ... $\mathbf{d}(\mathbf{f}_1(\mathbf{x}) + \mathbf{f}_2(\mathbf{x})) / \mathbf{d}\mathbf{x} = \mathbf{df}_1(\mathbf{x}) / \mathbf{d}\mathbf{x} + \mathbf{df}_2(\mathbf{x}) / \mathbf{d}\mathbf{x}$







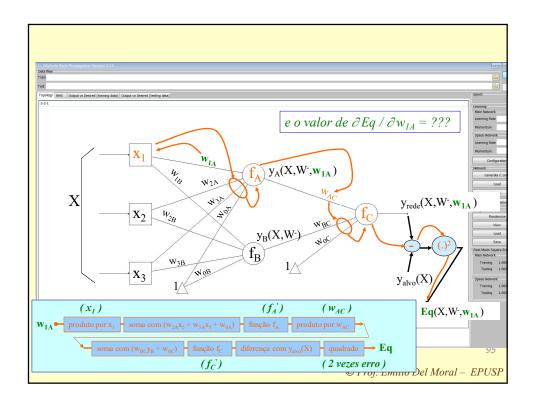


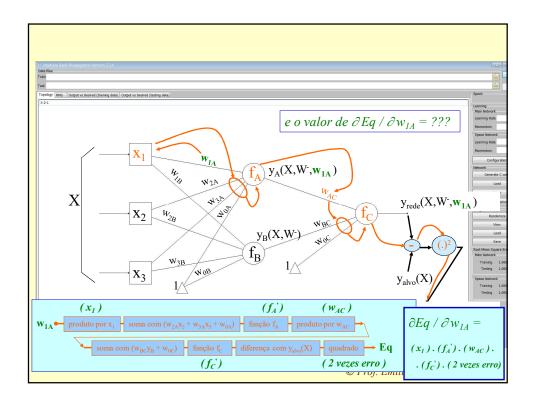
Note que aqui temos uma cadeia com muitos estágios que levam da váriável W_{1A} , à variável Eq^{μ} , e para a qual podemos calcular a derivada da saída (Eq^{μ}) com relação à entrada (W_{1A}) aplicando de forma repetida a seguinte propriedade simples e sua velha conhecida ...

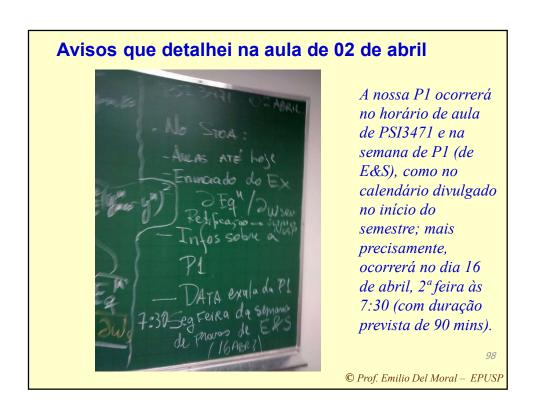
$$d(f_1(f_2(x)) / dx = df_1(x)/df_2 \cdot df_2(x)/dx$$

..., ou seja, calculando isoladamente o valor da derivada para cada estágio da cadeia, e finalizando o cálculo de derivada de ponta a ponta nessa cadeia toda através do produto dos diversos valores de cada estágio.

94









Lembretes

- Na maioria dos slides anteriores, onde aparece X, leia-se
 X^μ, não incluído para não complicar demais os desenhos
- ... similarmente, onde aparece y_{alvo} , leia-se y_{alvo} $^{\mu}$. Idem para os Eq, leia-se Eq $^{\mu}$
- Nos itens de cadeia de derivadas (f_A) e (f_C) , atenção para os valores dos argumentos, que devem ser os mesmos de f_A e f_C na cadeia original que leva w_{IA} a Eq.
- ... lembrando ... na cadeia original tinhamos ...
 - para f_C : $f_C(w_{AC} \cdot f_A(w_{1A}.x_1 + w_{2A}.x_2 + w_{0A}) + w_{BC} \cdot f_B(w_{1B}.x_1 + w_{2B}.x_2 + w_{0B}) + w_{0C})$ - para f_A : $f_A(w_{1A}.x_1 + w_{2A}.x_2 + w_{0A})$
- Similarmente, para o bloco "quadrado", cuja derivada é a função "2 vezes erro", o argumento é [y_{rede}(X,W) y_{alvo}(X)]

101

© Prof. Emilio Del Moral – EPUSP

Lembretes

- O mesmo que foi feito para w_{IA} deve ser feito agora para os demais 10 pesos: w_{2A} , w_{3A} , $w_{\theta A}$, w_{1B} , w_{2B} , w_{3B} , $w_{\theta B}$, w_{AC} , w_{BC} , e $w_{\theta C}$!
- Assim compomos um gradiente de 11 dimensões, com as derivadas de Eq^μ com relação aos 11 diferentes pesos w: Grad_w (Eq^μ)
- Essas 11 fórmulas devem ser aplicadas repetidamente aos M exemplares numéricos de X^μ e y_{alvo}^μ, calculando M gradientes!
- Com eles, se obtém o gradiente médio dos M pares empíricos: Grad_w (Eqm) = [Σ_{ιι} Grad_w (Eq^μ)] / M
- Esse gradiente médio é a Bussola do Gradiente!

102

Método do Gradiente Aplicado aos nossos MLPs: a partir de um W#0, temos aproximações sucessivas ao Eqm mínimo, por repetidos pequenos passos DeltaW, sempre contrários ao gradiente ...

- "Chute" um W inicial para o "Wcorrente", ou "W melhor até agora"
- Em loop até obter Egm zero, ou baixo o suficiente, ou estável:
 - Determine o vetor gradiente do Eqm, nesse espaço de Ws
 - Em loop varrendo todos os M exemplos (X^μ;y^μ),
 - Calcule o gradiente de Eq $^\mu$ associado a um exemplo μ , e vá varrendo μ e somando os gradientes de cada Eq $^\mu$, para compor o vetor gradiente de Eqm, assim que sair deste loop em μ ;
 - Cada cálculo como esse, envolve primeiro calcular os argumentos de cada tangente hiperbólica e depois usar esses argumentos na regra da cadeia das derivadas necessárias
 - Dê um passo Delta ΔW nesse espaço, com direção e magnitude dados por −η*vetor gradiente médio para os M Exemplos (X^μ;y^μ) de treino

103



O Ciclo completo da modelagem:

- 0) Formalização do problema, mapeamento quantitativo em um modelo neural inicial e ... 0b) coleta de pares empíricos (X,y)
- 1) Fase de TREINO da RNA (MLP): com conhecimento dos X e dos y, que são ambos usados na calibração do modelo
- 2) Fase de TESTE / Caracterização da qualidade da RNA para generalizar: temos <u>novos</u> pares X e y, com y guardado "na gaveta", usado apenas para avaliação, não para re-calibração. É como um ensaio de uso final do modelo, com possibilidade de medir a sua qualidade com o y que foi guardado na gaveta.

[Fase de refinamentos da RNA, dados e modelo, em ciclos, desde 0]

3) Fase de USO FINAL da RNA, com y efetivamente não conhecido, e estimado com conhecimento dos X + uso do modelo calibrado.

.... Diferenças e semelhanças entre 1, 2 e 3

107

